

CUNY SPS DATA 621 - CTG5 - HW3

Betsy Rosalen, Gabrielle Bartomeo, Jeremy O'Brien, Lidiia Tronina, Rose Koh

April 10th, 2019

Contents

1	DATA EXPLORATION	2
1.1	Summary Statistics	3
1.2	Shape of Predictor Distributions	3
1.3	Outliers	4
1.4	Missing Values	5
1.5	Linearity	5
2	DATA PREPARATION	7
2.1	Missing Values and NA Imputation	7
2.2	Dealing with outliers, leverage, and influence points	7
2.3	Correlation	9
2.4	Feature Engineering	10
3	BUILD MODELS	11
3.1	Model 1 - Base Model	11
3.2	Model 2 - Log Transform Skewed Predictors and Automated Selection Tools	13
3.3	Model 3 - Examine all possible interactions	16
3.4	Model 4 - Segmented/Piecewise Regression	17
4	SELECT MODELS	22
4.1	Pseudo R2	22
4.2	Summary diagnostic plots	23
5	Appendix	24

1 DATA EXPLORATION

Relocating to a new city or state can be very stressful. In addition to the stress of packing and moving, you may also be nervous about moving to an unfamiliar area. To better understand their new community, some new residents or people interested in moving to a new city choose to review crime statistics in and around their neighborhood. Crime rate may also influence where people choose to live, raise their families and run their businesses; many potential new residents steer clear of cities with higher than average crime rates.

Data was collected in order to predict whether the neighborhood will be at risk for high crime levels. For each neighborhood the response variable, **target**, represents whether the crime rate is above the median crime rate or not. In addition to that 13 predictor variables were collected representing each neighborhood's: proportion of large lots, non-retail business acres, whether or not it borders the Charles River, nitrogen oxides concentration, average number of rooms per dwelling, proportion of owner-occupied units, distances to five Boston employment centers, accessibility to radial highways, property tax rate, pupil-teacher ratio, proportion of African Americans, percent lower status, and median value of homes. The evaluation data contains the same 13 predictor variables and no target variable so it will be impossible to check the accuracy of our predictions from the testing data.

Table 1: Data Dictionary

VARIABLE	DEFINITION	TYPE
target	whether the crime rate is above the median crime rate (1) or not (0)	response
zn	proportion of residential land zoned for large lots (over 25000 square feet)	predictor
indus	proportion of non-retail business acres per suburb	predictor
chas	a dummy var. for whether the suburb borders the Charles River (1) or not (0)	predictor
nox	nitrogen oxides concentration (parts per 10 million)	predictor
rm	average number of rooms per dwelling	predictor
age	proportion of owner-occupied units built prior to 1940	predictor
dis	weighted mean of distances to five Boston employment centers	predictor
rad	index of accessibility to radial highways	predictor
tax	full-value property-tax rate per \$10,000	predictor
ptratio	pupil-teacher ratio by town	predictor
black	$1000(B_k - 0.63)^2$ where B_k is the proportion of blacks by town	predictor
lstat	lower status of the population (percent)	predictor
medv	median value of owner-occupied homes in \$1000s	predictor

1.1 Summary Statistics

Table 2: Summary statistics

	n	min	mean	median	max	sd
zn	466	0.0000	11.5772532	0.00000	100.0000	23.3646511
indus	466	0.4600	11.1050215	9.69000	27.7400	6.8458549
chas	466	0.0000	0.0708155	0.00000	1.0000	0.2567920
nox	466	0.3890	0.5543105	0.53800	0.8710	0.1166667
rm	466	3.8630	6.2906738	6.21000	8.7800	0.7048513
age	466	2.9000	68.3675966	77.15000	100.0000	28.3213784
dis	466	1.1296	3.7956929	3.19095	12.1265	2.1069496
rad	466	1.0000	9.5300429	5.00000	24.0000	8.6859272
tax	466	187.0000	409.5021459	334.50000	711.0000	167.9000887
ptratio	466	12.6000	18.3984979	18.90000	22.0000	2.1968447
lstat	466	1.7300	12.6314592	11.35000	37.9700	7.1018907
medv	466	5.0000	22.5892704	21.20000	50.0000	9.2396814
target	466	0.0000	0.4914163	0.00000	1.0000	0.5004636

Looking at the Table. 1, we can see that **chas** and **target** are binary variables. 49% of our target variable is coded as 0's indicating that the crime rate is NOT above the median crime rate. There are potential outliers present in **zn**, **lstat**, **medv** and **dis**.

1.2 Shape of Predictor Distributions

Figure. 1 shows that the distribution of most of the variables seems skewed. There are some outliers in the right tail of **tax**, **rad**, **medv**, **lstat**, **dis** and left tail of **ptratio**.

Even more interestingly, for many of the predictor variables the shape of the distribution is significantly different depending on the value of the **target**. For example, **age** (proportion of owner-occupied units built prior to 1940) is highly left skewed for homes where the crime rate is above the median crime rate (**target** = 1), while for homes where the crime rate is not above the median (**target** = 0) the distribution is normal. This indicates that areas with a higher proportion of older structures are more likely to have a crime rate above the median, which is what we would expect. Other variables with similar differences are **dis**, **indus**, **lstat**, **nox**, **ptratio**, **rad**, and **tax**.

The variable **rad** has a clear separation at about a value of 5 with almost all homes with a **rad** value less than five being in the **target** group coded 0 and almost all homes with a **rad** value greater than five being in the **target** group coded 1. Possibly indicating that a transformation into a categorical dummy variable might be desirable. **indus** has a similar separation at a value of about 16, but not as strikingly.

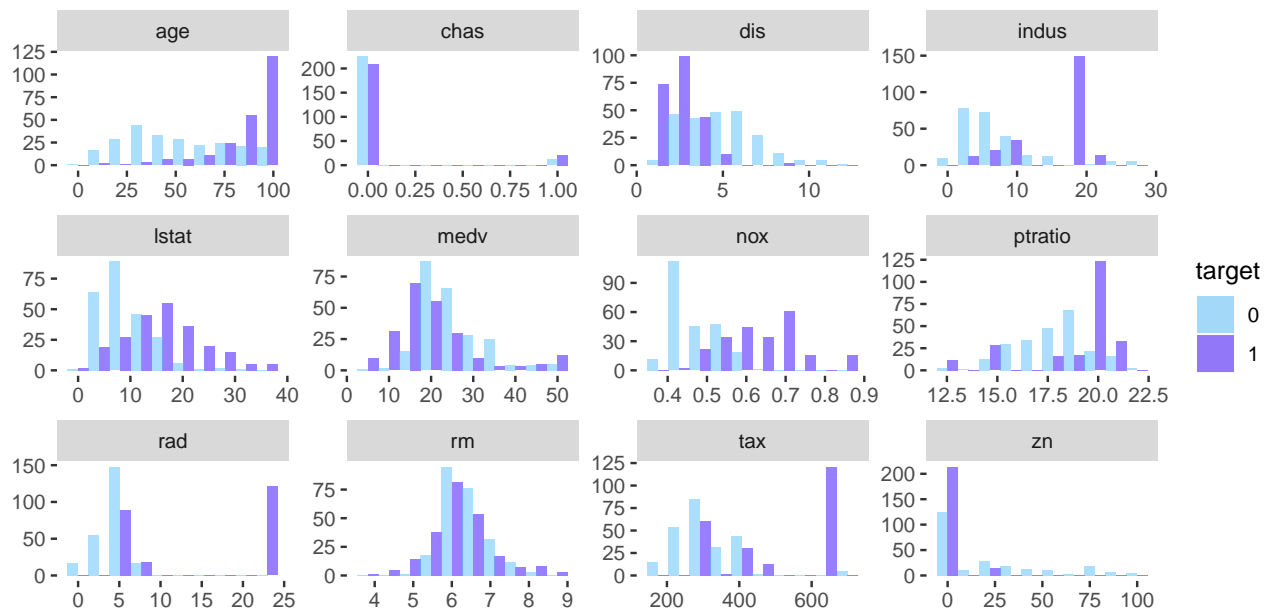


Figure 1: Data Distributions

1.3 Outliers

Figure. 2 shows that there are also a large number of outliers that need to be accounted for, most significantly in `zn` (proportion of residential land zoned for large lots [over 25000 square feet]) and `medv` (median value of owner-occupied homes in \$1000s) and less significantly in `lstat`, `dis` and `rm`. Since the `tax` variable has values which are very large compared to other variables in the dataset, it was scaled to fit the boxplot by dividing by 10.

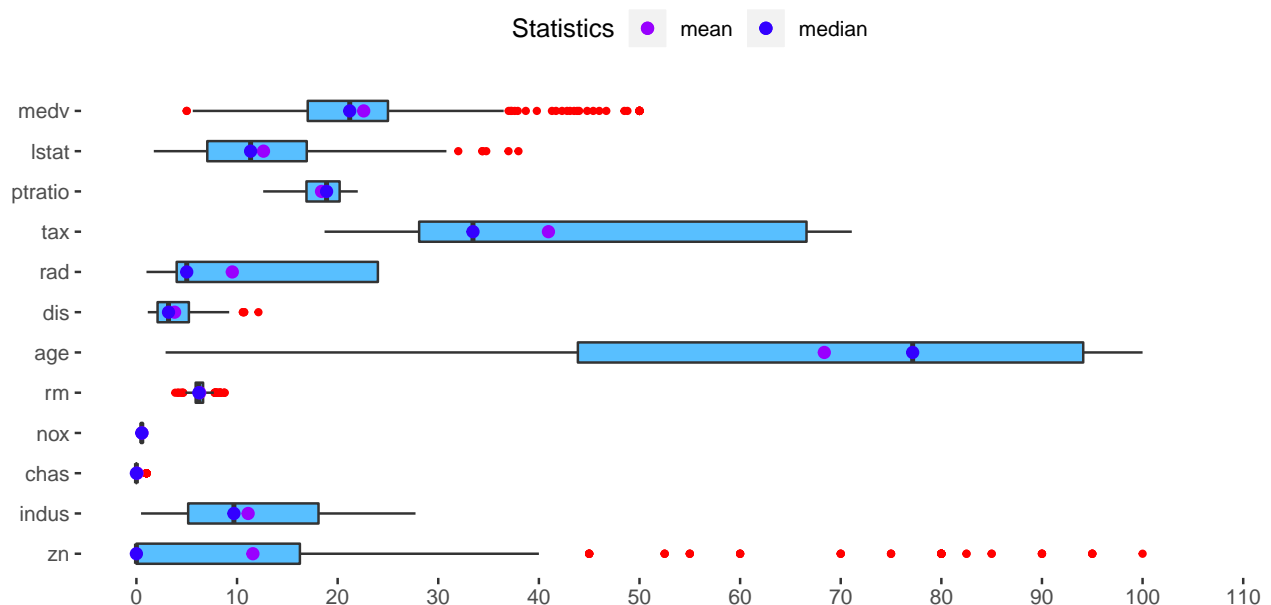


Figure 2: Boxplots highlighting many outliers in the data.

1.4 Missing Values

There are no missing values in any of our observations gathered across the thirteen predictor variables as can be seen in Figure. 3.

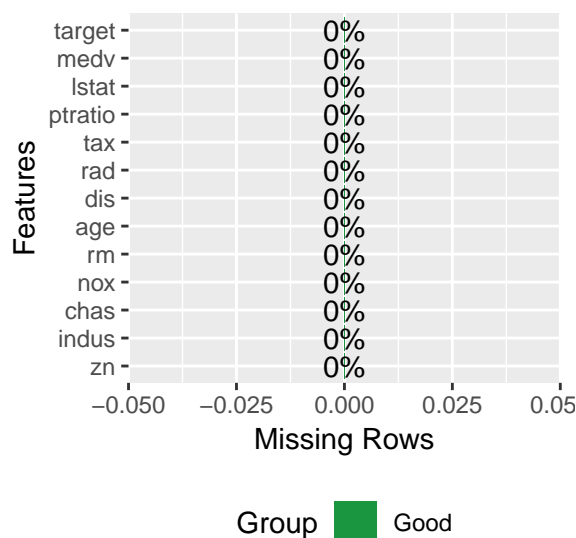


Figure 3: Missing values

1.5 Linearity

Each variable was plotted against the target variable in order to determine at a glance which had the most potential linearity before the dataset was modified.

As can be observed in Figure. 4, all of the predictor variables seem to have an impact on the target. With most of them having a positive impact indicating that the higher the predictor variable values are more likely to correspond to a target that is coded as 1 indicating the crime rate is above the median. The exceptions are:

1. **dis**, weighted mean of distances to five Boston employment centers
2. **medv**, median value of owner-occupied homes in \$1000s
3. **rm**, average number of rooms per dwelling
4. **zn**, proportion of residential land zoned for large lots (over 25000 square feet)
5. and possibly **chas**, a dummy var. for whether the suburb borders the Charles River (1) or not (0)

For these variables the distribution of predictor variable values is higher when the target is coded 0 for 'crime rate not above the median'.

We can also see that many of the predictor variables have very different variances for the two values of the target. This is especially true for **age**, **rad**, **tax**, and **zn** and less significantly for **dis** and **nox**. The presence of a large number of outliers **zn**, **medv**, **lstat**, **dis** and **rm** as noted above also becomes more apparent. The large number of outliers in **age** which were hidden in the previous plot also become visible here.

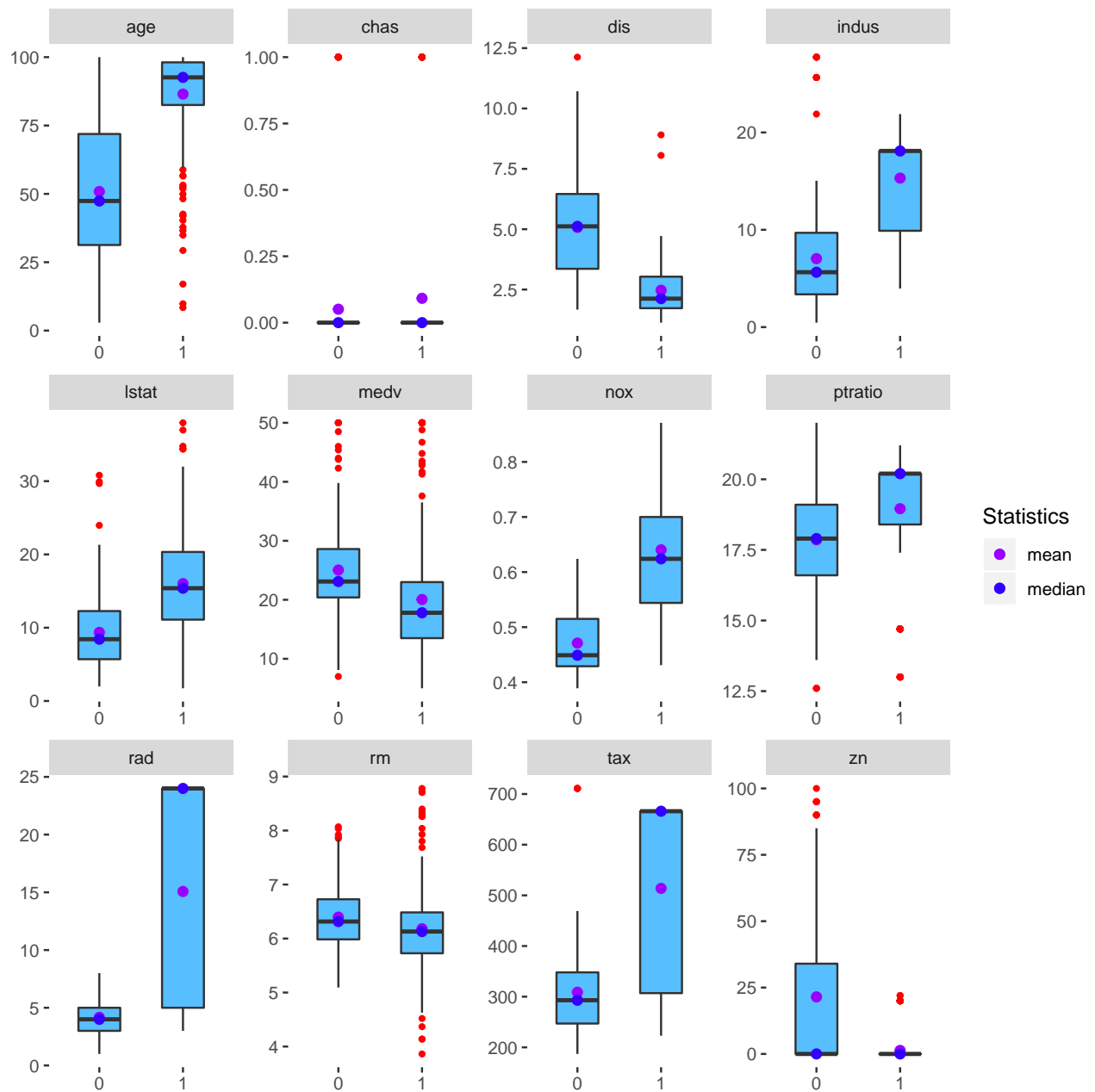


Figure 4: Linear relationships between each predictor and the target

2 DATA PREPARATION

2.1 Missing Values and NA Imputation

Given that (as noted above) the training dataset does not include any missing values, there's no need to make systematic corrections or imputations.

2.2 Dealing with outliers, leverage, and influence points

While logistic regression can be more robust to leverage points (explanatory variable values, which are distant on the x-axis), outliers (response variable values, which are distant on the y-axis) can exert influence which affects the curve and accuracy of target predictions.

- **dis**, **tax** (property tax rate per \$10k), and **medv** (median value of owner-occupied homes) see a few outliers and leverage points in both target classes
- **indus** (the non-retail business acreage proportion) and **lstat** (percent lower status population) both have outliers in the below-mean (0) class
- **prratio** (pupil-teacher ratio) fit is very impacted by density of low values in the above-mean class, making the linear relationship appear parabolic
- **rad** (highway access index) is influenced by a high-value concentration of locations distant from radial highways that fall in the above-mean class and is almost a perfect predictor for our target with almost all values at 5 and above being coded 1 indicating they are above the median crime rate.
- **rm** (average rooms per dwelling) sees a wider distribution of house size for the above-mean class; while **zn** (large-lot zoned land proportion) sees the opposite, with a concentration around a few non-residential land proportions for the above-mean class and a wide dispersion for the below-mean class

The figures below examine the linear relationships after a log transformation, which smoothes several relationships but still demonstrates visible influence for several other variables: **lstat**, **medv**, **prratio**, **rad**, **rm**, **tax**, and **zn**. We discuss further in the feature engineering section below.

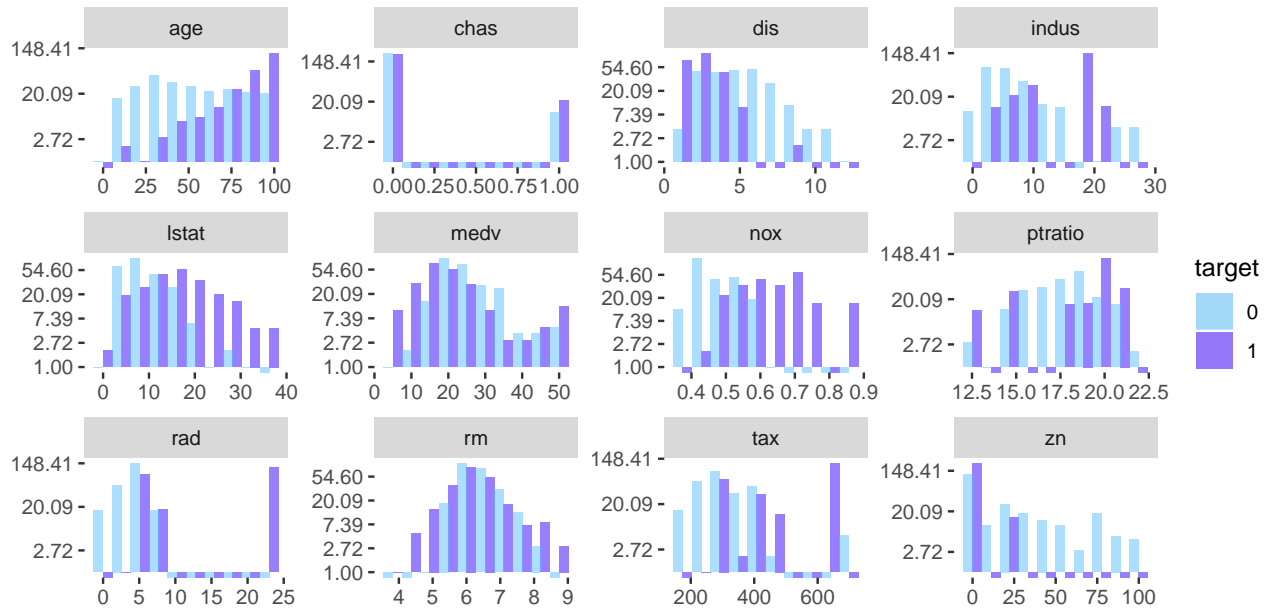


Figure 5: Natural log transformed predictor distributions

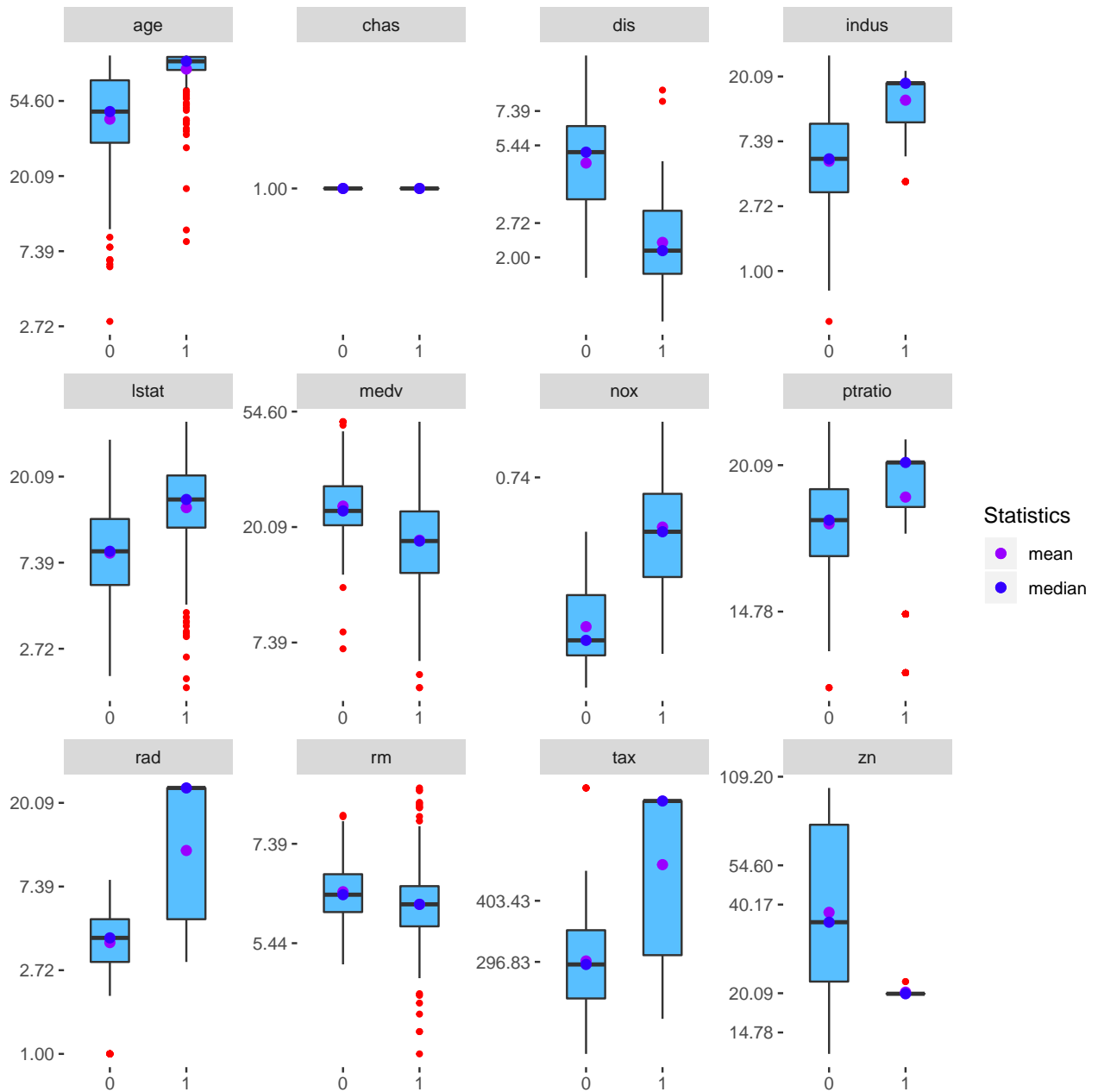


Figure 6: Relationships between natural log transformed predictors and the target

2.3 Correlation

An examination of correlation between the explanatory variables reveals the following:

- **indus** (non-retail business acre proportion) is positively correlated with **nox** (pollution concentration, $r = .76$) and **tax** (property tax rate per \$10k, $r = .73$) and is negatively correlated with **dis** (weighted mean distance to employment centers, $r = -.7$)
- **chas** (bordering Charles river) correlated with **nox** ($r = .97$) and **rm** (average rooms per dwelling, $r = .91$) and **age** (proportion of pre-1940 homes, $r = .79$); and is negatively correlated with **dis** ($r = -.97$)
- **medv** (median value of owner-occupied homes) is correlated with **rm** ($r = .71$); and is negatively correlated with **lstat** (percent lower status population, $r = -.74$)
- **age** is correlated with **nox** ($r = .74$); and is negatively correlated with **dis** ($r = -.75$)
- **rad** (highway access index) correlated with **tax** ($r = .91$)

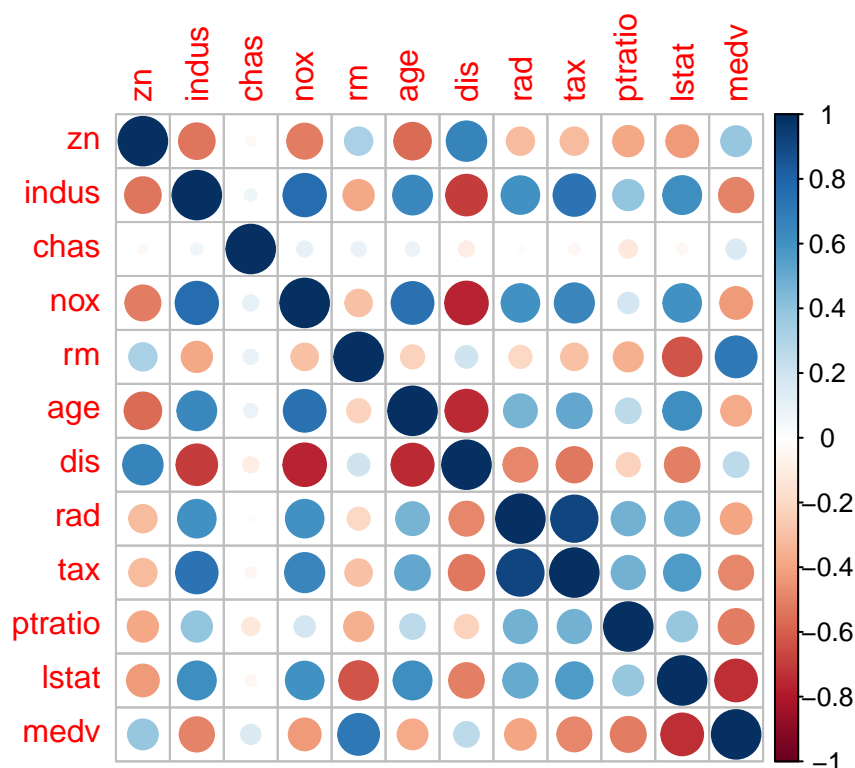


Table 3: Correlation between predictors

	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	lstat	medv
zn	1.00	-0.54	-0.04	-0.52	0.32	-0.57	0.66	-0.32	-0.32	-0.39	-0.43	0.38
indus	-0.54	1.00	0.06	0.76	-0.39	0.64	-0.70	0.60	0.73	0.39	0.61	-0.50
chas	-0.04	0.06	1.00	0.10	0.09	0.08	-0.10	-0.02	-0.05	-0.13	-0.05	0.16
nox	-0.52	0.76	0.10	1.00	-0.30	0.74	-0.77	0.60	0.65	0.18	0.60	-0.43
rm	0.32	-0.39	0.09	-0.30	1.00	-0.23	0.20	-0.21	-0.30	-0.36	-0.63	0.71
age	-0.57	0.64	0.08	0.74	-0.23	1.00	-0.75	0.46	0.51	0.26	0.61	-0.38
dis	0.66	-0.70	-0.10	-0.77	0.20	-0.75	1.00	-0.49	-0.53	-0.23	-0.51	0.26
rad	-0.32	0.60	-0.02	0.60	-0.21	0.46	-0.49	1.00	0.91	0.47	0.50	-0.40
tax	-0.32	0.73	-0.05	0.65	-0.30	0.51	-0.53	0.91	1.00	0.47	0.56	-0.49
ptratio	-0.39	0.39	-0.13	0.18	-0.36	0.26	-0.23	0.47	0.47	1.00	0.38	-0.52
lstat	-0.43	0.61	-0.05	0.60	-0.63	0.61	-0.51	0.50	0.56	0.38	1.00	-0.74
medv	0.38	-0.50	0.16	-0.43	0.71	-0.38	0.26	-0.40	-0.49	-0.52	-0.74	1.00

2.4 Feature Engineering

In ‘A Modern Approach to Regression with R’ (page 284), Sheather quotes Cook and Weisberg, suggesting that the best way to determine need for log transformation of skewed predictors is to include both the original and transformed variables in the logistic regression model in order assess their relative contributions directly and prune accordingly

Reexamining the histograms of the predictor distributions above reveals that:

- **age** is left-skewed
- **dis** is right-skewed, and **zn** is extremely so
- **nox** is right-skewed and platykurtic (thin-tailed)
- **rad** and **tax** seem to have normal distributions, with large numbers of outliers at particular levels
- **indus** and **ptratio** reveal peculiar skew, with incidences at particular high level, perhaps due to regulation or infrastructure requirements

We include log transforms of **age**, **dis**, **nox**, **rad**, **tax**, **indus**, and **ptratio** in the dataset for evaluation in models.

We also looked at interactions between variables in two ways; first by including all possible interactions in one model, and then by choosing only the most highly correlated variables to include as interactions in a subsequent model.

Finally, we also separated the observations with **rad** values at 5 and above vs. those with less than 5 to create a separate regression line in a segmented regression approach.

3 BUILD MODELS

3.1 Model 1 - Base Model

The First model is a binary logistic model including all the explanatory variables. The data is centered and scaled based on the mean and standard deviation of the variables. The residual deviance is 192.05, AIC is 218.05, Pseudo- R^2 is 0.83 or 0.70 depending on which method you use for calculation. We will consider this as the baseline for all models. Two of the variables **rm** and **medv** have VIF values above the usual cutoff at 5 indicating that collinearity may be causing problems in our model. We will address this in subsequent models.

Observations	466
Dependent variable	target
Type	Generalized linear model
Family	binomial
Link	logit

$\chi^2(12)$	453.83
Pseudo- R^2 (Cragg-Uhler)	0.83
Pseudo- R^2 (McFadden)	0.70
AIC	218.05
BIC	271.92

	Est.	S.E.	z val.	p	VIF
(Intercept)	-40.82	6.63	-6.15	0.00	NA
zn	-0.07	0.03	-1.90	0.06	1.82
indus	-0.06	0.05	-1.36	0.17	2.68
chas	0.91	0.76	1.21	0.23	1.24
nox	49.12	7.93	6.19	0.00	4.16
rm	-0.59	0.72	-0.81	0.42	5.81
age	0.03	0.01	2.47	0.01	2.57
dis	0.74	0.23	3.21	0.00	3.89
rad	0.67	0.16	4.08	0.00	1.94
tax	-0.01	0.00	-2.09	0.04	2.14
ptratio	0.40	0.13	3.18	0.00	2.28
lstat	0.05	0.05	0.85	0.40	2.64
medv	0.18	0.07	2.65	0.01	8.12

Standard errors: MLE

$$\hat{y} = 2.33 - 1.54\text{zn} - 0.44\text{indus} + 0.23\text{chas} + 5.73\text{nox} - 0.41\text{rm} + 0.97\text{age} \\ + 1.55\text{dis} + 5.79\text{rad} - 1.04\text{tax} + 0.88\text{ptratio} + 0.33\text{lstat} + 1.67\text{medv}$$

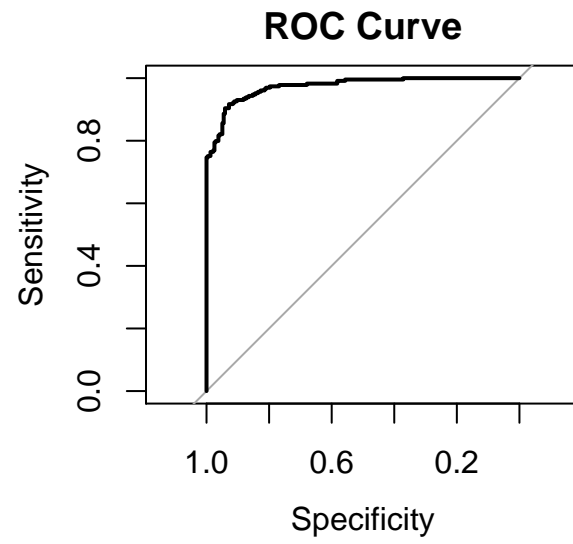


Figure 7: Model 1 ROC Curve

Table 4: Area Under the Curve

x
0.9737623

3.2 Model 2 - Log Transform Skewed Predictors and Automated Selection Tools

The second model is a binary logistic model including all the explanatory variables plus log transformations of our skewed variables `age`, `dis`, `nox`, `rad`, `tax`, `indus`, and `ptratio` as recommended by Sheather in ‘A Modern Approach to Regression with R’. We see considerably improved statistics below with higher Pseudo- R^2 values for both the (Craig-Uhler) and (McFadden) calculations and lower AIC and BIC values. We see some very large VIF numbers though indicating that we have too many correlated variables in our model before refinement. So we used the step function to refine our model next.

Observations	466
Dependent variable	target
Type	Generalized linear model
Family	binomial
Link	logit

$\chi^2(19)$	512.55
Pseudo- R^2 (Cragg-Uhler)	0.89
Pseudo- R^2 (McFadden)	0.79
AIC	173.33
BIC	256.21

	Est.	S.E.	z val.	p	VIF
(Intercept)	-322.87	145.89	-2.21	0.03	NA
zn	-0.05	0.06	-0.93	0.35	4.21
indus	0.38	0.28	1.33	0.18	37.50
chas	-0.51	0.97	-0.53	0.60	1.73
nox	230.71	124.56	1.85	0.06	827.63
rm	-1.73	1.01	-1.71	0.09	5.23
age	0.10	0.04	2.67	0.01	18.66
dis	-2.74	1.04	-2.63	0.01	54.37
rad	1.54	0.39	3.92	0.00	20.09
tax	-0.20	0.05	-3.59	0.00	268.34
ptratio	5.31	2.25	2.36	0.02	412.82
lstat	0.03	0.06	0.52	0.61	2.37
medv	0.29	0.11	2.74	0.01	6.92
log(age)	-2.97	1.47	-2.01	0.04	17.77
log(dis)	13.41	4.45	3.02	0.00	64.05
log(nox)	-98.55	64.18	-1.54	0.12	812.46
log(rad)	-1.32	1.52	-0.87	0.39	8.16
log(tax)	58.72	16.96	3.46	0.00	203.30
log(indus)	-1.55	2.86	-0.54	0.59	31.61
log(ptratio)	-81.99	39.96	-2.05	0.04	417.11

Standard errors: MLE

3.2.1 Refining with the step function and backward elimination

Using the step function to refine this model leaves us with the following model plus `log(nox)`. Since `log(nox)` had such a high p-value we removed it from the model, which didn’t change the model’s performance much.

Observations	466
Dependent variable	target
Type	Generalized linear model
Family	binomial
Link	logit

$\chi^2(13)$	508.74
Pseudo- R^2 (Cragg-Uhler)	0.89
Pseudo- R^2 (McFadden)	0.79
AIC	165.14
BIC	223.16

	Est.	S.E.	z val.	p	VIF
(Intercept)	-111.79	78.84	-1.42	0.16	NA
indus	0.28	0.09	2.95	0.00	4.67
nox	44.58	9.94	4.48	0.00	5.31
rm	-1.54	0.81	-1.90	0.06	3.79
age	0.11	0.03	3.46	0.00	13.89
dis	-2.30	0.75	-3.08	0.00	28.91
rad	1.32	0.25	5.29	0.00	9.39
tax	-0.18	0.04	-4.18	0.00	189.78
ptratio	5.87	2.11	2.78	0.01	395.34
medv	0.25	0.10	2.58	0.01	6.07
log(age)	-4.01	1.27	-3.16	0.00	12.60
log(dis)	11.12	3.34	3.33	0.00	38.27
log(tax)	52.55	12.50	4.20	0.00	127.33
log(ptratio)	-91.86	37.02	-2.48	0.01	393.71

Standard errors: MLE

Although our Pseudo- R^2 values didn't change at all we see an improvement in both the AIC and BIC numbers in a much simpler and reduced model which is preferred for simplicity sake.

The equation for the simplified second model is:

$$\begin{aligned}\hat{y} = & -111.79 + 0.28\text{indus} + 44.58\text{nox} - 1.54\text{rm} + 0.11\text{age} \\ & - 2.30\text{dis} + 1.32\text{rad} - 0.18\text{tax} + 5.87\text{ptratio} + 0.24\text{medv} \\ & - 4.01\log(\text{age}) + 11.12\log(\text{dis}) + 52.55\log(\text{tax}) - 91.86\log(\text{ptratio})\end{aligned}$$

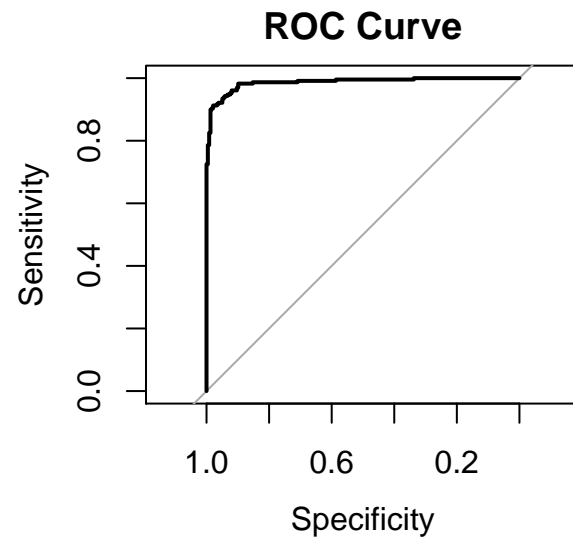


Figure 8: Model 2 ROC Curve

Table 5: Area Under the Curve

x
0.9866232

3.3 Model 3 - Examine all possible interactions

For Model 3 we started with a strategy suggested by Faraway in ‘Extending the Linear Model with R’ that adds all possible interactions between the predictor variables in addition to the full set of predictors then uses the step function to remove unnecessary variables or interactions. The resulting model is on the following page.

Observations	466
Dependent variable	target
Type	Generalized linear model
Family	binomial
Link	logit

$\chi^2(35)$	645.88
Pseudo- R^2 (Cragg-Uhler)	1.00
Pseudo- R^2 (McFadden)	1.00
AIC	72.00
BIC	221.19

Although we got the best AIC (72) and residual deviance, AIC, BIC and Pseudo- R^2 numbers with this model, the coefficients are all insanely large and so are the VIF’s and the p-values with not a single one showing any significance. This is a good example of an extremely over-fitted model. It probably models our training data perfectly, but would perform very poorly at predicting our test data.

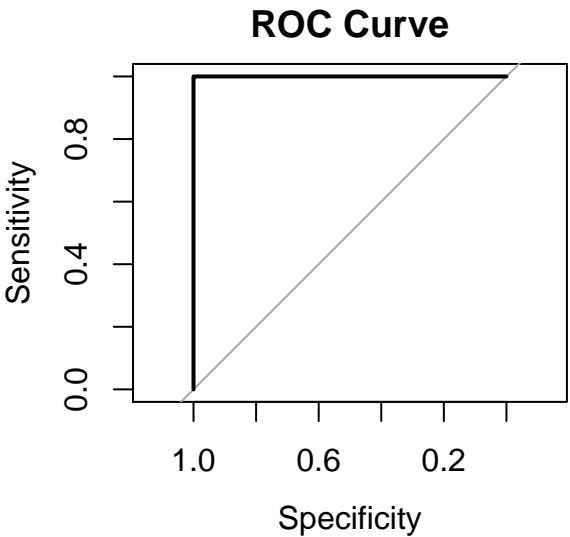


Figure 9: Model 3 ROC Curve

Table 6: Area Under the Curve

$\frac{x}{1}$

	Est.	S.E.	z val.	p	VIF
(Intercept)	3762.16	758565.87	0.00	1.00	NA
zn	330.82	10175.45	0.03	0.97	176034.62
indus	-825.21	18300.11	-0.05	0.96	162947.92
chas	560.68	30141.28	0.02	0.99	284.48
nox	61117.01	1331465.34	0.05	0.96	100549.64
rm	-525.21	17382.15	-0.03	0.98	573.08
age	-85.23	2926.21	-0.03	0.98	94911.42
dis	-4914.20	142609.93	-0.03	0.97	1077253.24
rad	148.06	2865.92	0.05	0.96	343.86
tax	-74.52	1584.62	-0.05	0.96	167240.72
ptratio	1110.07	25192.51	0.04	0.96	23276.20
lstat	734.72	14060.74	0.05	0.96	82138.67
medv	-923.70	20931.95	-0.04	0.96	127353.16
zn:age	-1.62	50.60	-0.03	0.97	3056.12
zn:tax	-0.33	12.64	-0.03	0.98	32879.63
zn:ptratio	-12.39	360.39	-0.03	0.97	64139.86
zn:lstat	6.00	235.87	0.03	0.98	4576.73
indus:chas	-80.51	2326.62	-0.03	0.97	398.60
indus:rad	46.08	1174.42	0.04	0.97	17861.28
indus:ptratio	42.02	935.74	0.04	0.96	193511.99
indus:medv	-4.73	207.04	-0.02	0.98	3854.74
nox:age	-122.29	2696.27	-0.05	0.96	30642.32
nox:tax	42.39	1159.81	0.04	0.97	58487.44
nox:ptratio	-3981.06	71996.82	-0.06	0.96	188499.05
nox:lstat	-846.33	16709.05	-0.05	0.96	51451.70
nox:medv	1192.12	32574.05	0.04	0.97	71744.83
rm:age	6.34	188.74	0.03	0.97	14622.04
age:tax	0.15	3.44	0.04	0.96	28225.40
age:ptratio	3.57	95.85	0.04	0.97	46950.20
dis:tax	6.29	148.47	0.04	0.97	102378.44
dis:ptratio	140.47	5034.84	0.03	0.98	465018.06
dis:lstat	-29.26	792.79	-0.04	0.97	3794.69
dis:medv	35.58	731.58	0.05	0.96	18612.24
rad:tax	-0.93	21.80	-0.04	0.97	6275.17
tax:medv	0.98	17.59	0.06	0.96	14668.14
lstat:medv	-8.43	160.37	-0.05	0.96	2780.97

Standard errors: MLE

3.4 Model 4 - Segmented/Piecewise Regression

Model 4 uses a ‘segmented’ or ‘piecewise’ approach based on the evidence we saw in the predictor distribution plots that indicated that splitting **rad** at a value of 5 would separate most 0’s from 1’s in our target. We also included the log transformed variables that remained at the end of the model 2 selection process and interaction terms for the variables that showed the greatest correlation in our correlation plot which we will later narrow down using the step function.

Clearly separating the data based on splitting the **rad** variable at a value of 5 showed a great improvement in the model with much better statistics all around. We still have some very high VIF numbers but those should be reduced after we refine the model by removing some of the variables. The two **rad** predictors for

Observations	466
Dependent variable	target
Type	Generalized linear model
Family	binomial
Link	logit

$\chi^2(24)$	557.46
Pseudo-R ² (Cragg-Uhler)	0.93
Pseudo-R ² (McFadden)	0.86
AIC	138.42
BIC	242.02

	Est.	S.E.	z val.	p	VIF
(Intercept)	-110.08	190.57	-0.58	0.56	NA
less_than_five(rad)	2.89	0.63	4.56	0.00	19.32
five_and_over(rad)	1.82	0.42	4.33	0.00	23.27
zn	-0.14	0.07	-1.96	0.05	7.33
indus	0.97	1.92	0.51	0.61	875.49
chas	-0.88	1.22	-0.72	0.47	1.88
nox	172.22	90.91	1.89	0.06	317.82
rm	-3.46	2.09	-1.65	0.10	16.29
age	0.58	0.35	1.68	0.09	1051.99
dis	6.67	5.07	1.32	0.19	904.27
tax	-0.30	0.10	-2.97	0.00	404.74
ptratio	11.55	5.24	2.20	0.03	989.69
lstat	0.03	0.09	0.33	0.74	2.99
medv	-0.23	0.50	-0.46	0.65	98.23
log(age)	-2.62	2.14	-1.23	0.22	25.21
log(dis)	29.14	12.81	2.27	0.02	321.69
log(tax)	84.00	24.86	3.38	0.00	214.92
log(ptratio)	-203.38	93.97	-2.16	0.03	1008.87
indus:nox	-1.99	3.63	-0.55	0.58	1337.05
indus:dis	-0.21	0.16	-1.34	0.18	42.50
indus:tax	0.00	0.00	0.83	0.41	404.54
nox:age	-0.76	0.56	-1.37	0.17	954.60
nox:dis	-20.94	13.57	-1.54	0.12	972.89
rm:medv	0.07	0.07	1.04	0.30	141.90
age:dis	-0.03	0.02	-1.63	0.10	29.41

Standard errors: MLE

values 5 and above or less than 5 are by far the two most significant predictors in our model. There are still a lot of variables with very low significance though and the model could definitely use some refinement.

3.4.1 Backward Elimination vs. Forward Selection

We tried using the step function for both forward selection and backward elimination and found that the backward elimination process resulted in the better model. There were still some variables that appeared to be adding little value to the model however, so we removed `dis`, then `medv`, then the intercept resulting in the final model below.

While we have much improved statistics over all, we still have some extremely high VIF numbers indicting a lot of collinearity in our model. This is to be expected since we still have some original and log transformed versions of the same variable, or original and interaction terms that include that variable left in the same model.

Observations	466
Dependent variable	target
Type	Generalized linear model
Family	binomial
Link	logit

$\chi^2(14)$	551.48
Pseudo-R ² (Cragg-Uhler)	0.93
Pseudo-R ² (McFadden)	0.85
AIC	124.39
BIC	186.56

	Est.	S.E.	z val.	p	VIF
less_than_five(rad)	3.01	0.51	5.85	0.00	NA
five_and_over(rad)	1.80	0.34	5.30	0.00	24.57
zn	-0.15	0.05	-2.84	0.00	32.45
indus	0.53	0.27	1.92	0.05	4.48
nox	32.55	11.57	2.81	0.00	115.07
rm	-3.59	1.47	-2.45	0.01	504.54
age	0.10	0.05	2.24	0.03	1114.76
tax	-0.22	0.06	-3.91	0.00	138.45
ptratio	10.87	2.62	4.15	0.00	4707.26
log(dis)	11.52	4.10	2.81	0.00	32311.85
log(tax)	69.02	18.06	3.82	0.00	404.31
log(ptratio)	-190.90	46.28	-4.12	0.00	147550.87
indus:dis	-0.18	0.09	-2.00	0.05	247646.87
rm:medv	0.05	0.02	2.76	0.01	135.79
age:dis	-0.02	0.01	-1.82	0.07	90.42

Standard errors: MLE

The equation for the fourth model is:

$$\begin{aligned}\hat{y} = & 0 + 3.01\text{rad (less than 5)} + 1.80\text{rad (five and over)} - 0.15\text{zn} \\ & + 0.53\text{indus} + 32.55\text{nox} - 3.59\text{rm} + 0.10\text{age} - 0.22\text{tax} + 10.87\text{ptratio} + 11.52\log(\text{dis}) \\ & + 69.02\log(\text{tax}) - 190.90\log(\text{ptratio}) - 0.18\text{indus:dis} + 0.04\text{rm:medv} + -0.02\text{age:dis}\end{aligned}$$

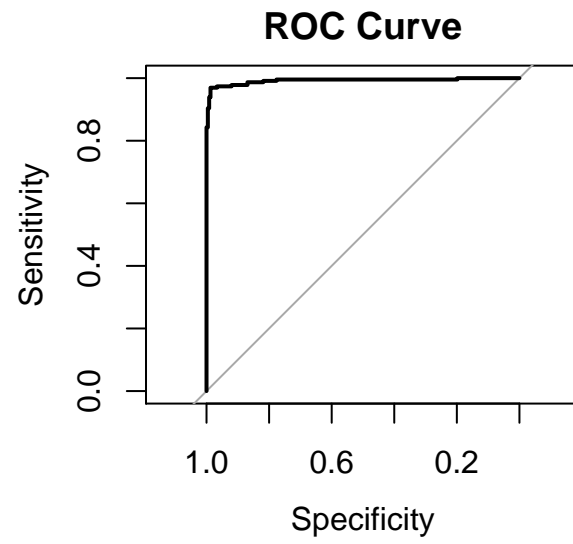
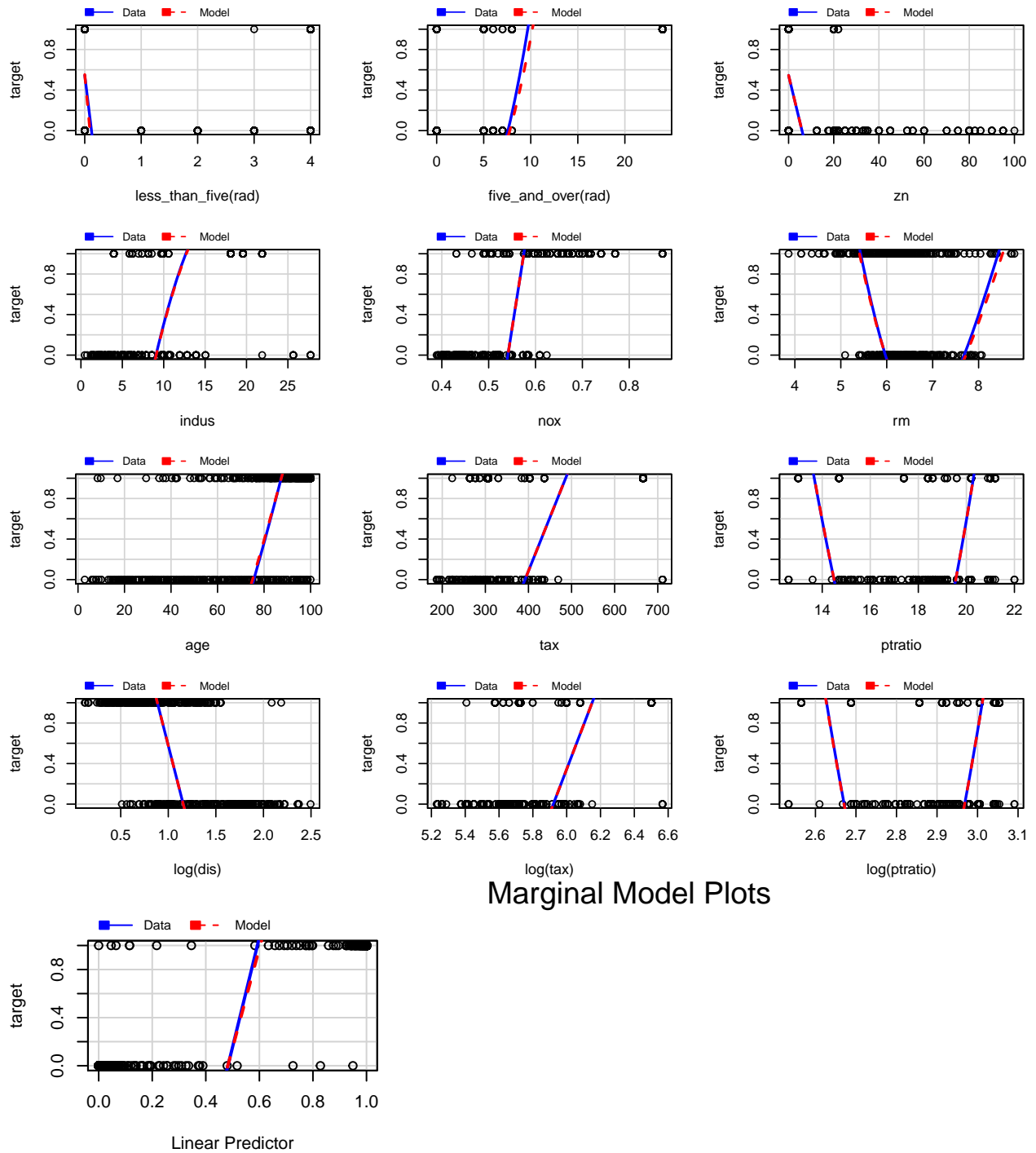


Figure 10: Model 4 ROC Curve

Table 7: Area Under the Curve

x
0.9921508

3.4.2 Marginal Model Plots



4 SELECT MODELS

Table 8: Confusion Matrix Summary Statistics

	Sensitivity	Specificity	Precision	Recall	F1
Model.1	0.9282700	0.9039301	0.9090909	0.9282700	0.9185804
Model.2	0.9662447	0.9126638	0.9196787	0.9662447	0.9423868
Model.3	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
Model.4	0.9873418	0.9650655	0.9669421	0.9873418	0.9770355

The four models were explored in order to determine the best way to determine whether or not a neighborhood's crime rate was above or below the median crime rate. It has been established that the most efficient model was the fourth model, with the first model being somewhat efficient, and the third over-fitted model being least efficient.

4.1 Pseudo R2

There is no R^2 for logistic regression to further evaluate, however, there is an alternative called *pseudoR²* terms that can be used for evaluation.

Table 9: Pseudo R2

	llh	llhNull	G2	McFadden	r2ML	r2CU
Model.1	-96.023459	-322.9379	453.8289	0.7026566	0.6223856	0.8299290
Model.2	-68.568641	-322.9379	508.7385	0.7876724	0.6643592	0.8858993
Model.3	-0.000007	-322.9379	645.8758	1.0000000	0.7499263	1.0000000
Model.4	-47.196447	-322.9379	551.4829	0.8538529	0.6937764	0.9251261

All of these measures and especially McFadden support the anova test's conclusion that model 4 is our most efficient and effective model for predicting whether a neighborhood will be at risk for a higher than median crime rate.

4.2 Summary diagnostic plots

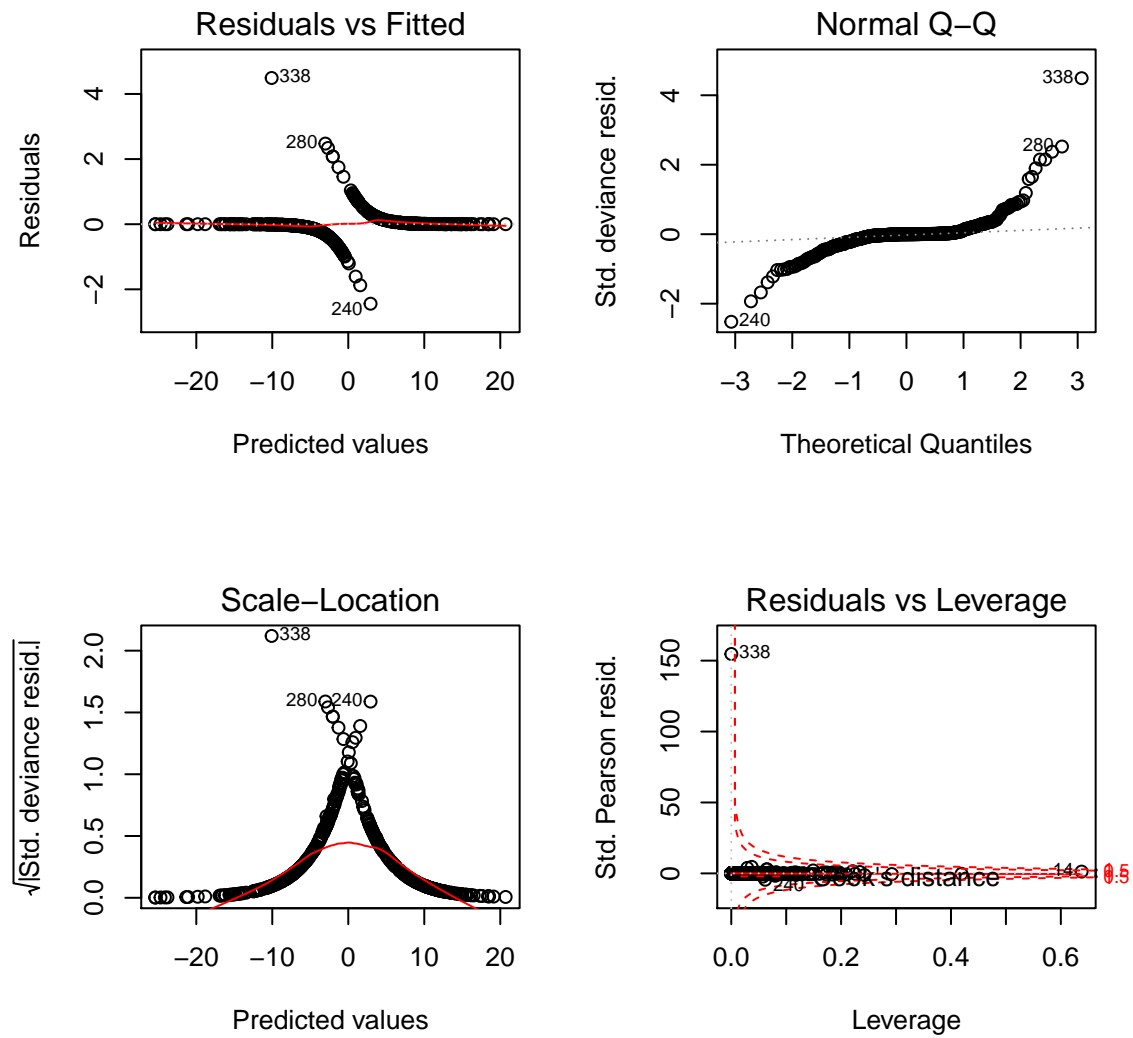


Figure 11: Model 4 Summary diagnostic plots

5 Appendix

The appendix is available as script.R file in `project3_crime` folder.

https://github.com/betsyrosalen/DATA_621_Business_Analyt_and_Data_Mining

```
# Load libs
```

```
if (!require('car')) (install.packages('car'))
if (!require('caret')) (install.packages('caret'))
if (!require('corrplot')) (install.packages('corrplot'))
if (!require('data.table')) (install.packages('data.table'))
if (!require('DataExplorer')) (install.packages('DataExplorer'))
if (!require('faraway')) (install.packages('faraway'))
if (!require('gridExtra')) (install.packages('gridExtra'))
if (!require('jtools')) (install.packages('jtools'))
if (!require('kableExtra')) (install.packages('kableExtra'))
if (!require('MASS')) (install.packages('MASS'))
if (!require('psych')) (install.packages('psych'))
if (!require('pROC')) (install.packages('pROC'))
if (!require('pscl')) (install.packages('pscl'))
if (!require('tidyverse')) (install.packages('tidyverse'))
```

```
# load data
```

```
train <- read.csv ('https://raw.githubusercontent.com/silverrainb/data621proj3/master/crime-training-da
```

```
test <- read.csv('https://raw.githubusercontent.com/silverrainb/data621proj3/master/crime-evaluation-da
```

```
variable_descriptions <- rbind(
```

```
  c('target','whether the crime rate is above the median crime rate (1) or not (0)','response'),
  c('zn','proportion of residential land zoned for large lots (over 25000 square feet) ','predictor'),
  c('indus','proportion of non-retail business acres per suburb','predictor'),
  c('chas','a dummy var. for whether the suburb borders the Charles River (1) or not (0)','predictor'),
  c('nox','nitrogen oxides concentration (parts per 10 million)','predictor'),
  c('rm','average number of rooms per dwelling','predictor' ),
  c('age','proportion of owner-occupied units built prior to 1940','predictor'),
  c('dis','weighted mean of distances to five Boston employment centers','predictor'),
  c('rad','index of accessibility to radial highways','predictor'),
  c('tax','full-value property-tax rate per $10,000','predictor'),
  c('ptratio','pupil-teacher ratio by town','predictor'),
  c('black',' $1000(B_k - 0.63)^2$  where  $B_k$  is the proportion of blacks by town','predictor'),
  c('lstat','lower status of the population (percent)','predictor'),
  c('medv','median value of owner-occupied homes in $1000s','predictor'))
```

```
colnames(variable_descriptions) <- c('VARIABLE','DEFINITION','TYPE')
```

```
# Summary Statistics
```

```
sum_stat <- describe(train)[,c(2,8,3,5,9,4)]
```

```
# Shape of Predictor Distributions
```

```
Hist_new <- train %>%
```

```
  gather(-target, key = "var", value = "val") %>%
  ggplot(aes(x = val, fill=factor(target))) +
  geom_histogram(position="dodge", bins=10, alpha=0.5) +
  facet_wrap(~ var, scales = "free") +
  scale_fill_manual("target",values = c("#58BFFF", "#3300FF")) +
  xlab("") +
```



```

model.2.raw <- glm(target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
  ptratio + lstat + medv + log(age) + log(dis) + log(nox) +
  log(rad) + log(tax) + log(indus) + log(ptratio),
  family = binomial,
  data = train)

mod2_summary_raw <- summ(model.2.raw, vifs = TRUE)

model.2.step <- step(model.2.raw, trace=FALSE)
mod2_summary_step <- summ(model.2.step, vifs = TRUE)

model.2 <- glm(target ~ indus + nox + rm + age + dis + rad + tax + ptratio + medv +
  log(age) + log(dis) + log(tax) + log(ptratio),
  family = binomial, data = train)

mod2_summary <- summary(model.2)
mod2_summary_a <- summ(model.2, vifs = TRUE)

#marg_mod_plot_2 <- mmeps(model.2, layout=c(5,4), key=NULL) # library car

### Model 2 Summary Statistics
pred.2.raw <- predict(model.2, newdata = train)
pred.2 <- as.factor(ifelse(pred.2.raw < .5, 0, 1))
mod2.conf.mat <- confusionMatrix(pred.2,
  as.factor(train$target), mode = "everything")

#=====#

## Model 5

#big_mod5 <- glm(target ~ (zn + indus + chas + nox + rm + age + dis + rad + tax +
#  ptratio + lstat + medv)^2, data = train, family = binomial)

#small_mod5 <- step(big_mod5, trace=FALSE)

# The above code is VERY computationally expensive
# Here's the result so it doesn't need to be run again.
model.5 <- glm(formula = target ~ zn + indus + chas + nox + rm + age + dis +
  rad + tax + ptratio + lstat + medv + zn:age + zn:tax + zn:ptratio +
  zn:lstat + indus:chas + indus:rad + indus:ptratio + indus:medv +
  nox:age + nox:tax + nox:ptratio + nox:lstat + nox:medv +
  rm:age + age:tax + age:ptratio + dis:tax + dis:ptratio +
  dis:lstat + dis:medv + rad:tax + tax:medv + lstat:medv,
  family = binomial,
  data = train)

mod5_summary <- summary(model.5)
mod5_summary_a <- summ(model.5, vifs = TRUE)

#resid_plot_5 <- residual.plots(model.5, exclude = 4, layout = c(2, 2)) # library car

#marg_mod_plot_5 <- mmeps(model.5, span = 3/4, layout = c(2, 2)) # library car

### Model 5 Summary Statistics

```

```

pred.5.raw <- predict(model.5, newdata = train)
pred.5 <- as.factor(ifelse(pred.5.raw < .5, 0, 1))
mod5.conf.mat <- confusionMatrix(pred.5, as.factor(train$target), mode = "everything")

#=====#
## Model 6

## Build the model
less_than_five <- function(x) ifelse(x < 5, x, 0)
five_and_over <- function(x) ifelse(x >= 5, x, 0)

model.6.raw <- glm(target ~ (less_than_five(rad) + five_and_over(rad)) + zn + indus + chas +
                    nox + rm + age + dis + tax + ptratio + lstat + medv + log(age) +
                    log(dis) + log(tax) + log(ptratio) + indus:nox + indus:dis +
                    indus:tax+ nox:age + nox:dis + rm:medv + dis:age,
                    family = binomial,
                    data = train)
mod6_summary_raw <- summ(model.6.raw, vifs = TRUE)
backward.mod <- step(model.6.raw, direction = "backward", trace=FALSE)
backward_sum <- summary(backward.mod)

#forward.mod <- step(model.6.raw, direction = "forward", trace=FALSE)
#forward_sum <- summary(forward.mod)

model.6 <- glm(target ~ less_than_five(rad) + five_and_over(rad) +
                zn + indus + nox + rm + age + tax + ptratio +
                log(dis) + log(tax) + log(ptratio) + indus:dis + rm:medv +
                age:dis + 0, family = binomial, data = train) # + 0 removes intercept

mod6_summary <- summary(model.6)
mod6_summary_a <- summ(model.6, vifs = TRUE)
### Model 6 Summary Statistics
pred.6.raw <- predict(model.6, newdata = train)
pred.6 <- as.factor(ifelse(pred.6.raw < .5, 0, 1))
mod6.conf.mat <- confusionMatrix(pred.6, as.factor(train$target), mode = "everything")

mod.6 <- train(target ~ less_than_five(rad) + five_and_over(rad) +
               zn + indus + nox + rm + age + tax + ptratio +
               log(dis) + log(tax) + log(ptratio) + indus:dis + rm:medv +
               age:dis + 0,
               family = binomial,
               data = train,
               method = 'glm') # + 0 removes intercept

#=====#

## Model Evaluations

eval_mods <- data.frame(mod1.conf.mat$byClass,
                        mod2.conf.mat$byClass,
                        mod5.conf.mat$byClass,
                        mod6.conf.mat$byClass) # add additional model stats

```

```
eval_mods <- data.frame(t(eval_mods))  
row.names(eval_mods) <- c("Model.1", "Model.2", "Model.3", "Model.4") # add additional models  
  
eval_mods <- dplyr::select(eval_mods, Sensitivity, Specificity, Precision, Recall, F1)  
  
# SELECT MODELS <<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<  
  
#Pseudo R2  
  
pseudo.r2 <- data.frame(pscl::pR2(model.1),  
                        pscl::pR2(model.2),  
                        pscl::pR2(model.5),  
                        pscl::pR2(model.6))  
  
pseudo.r2 <- data.frame(t(pseudo.r2))  
  
row.names(pseudo.r2) <- c("Model.1", "Model.2", "Model.3", "Model.4")
```